

BRIDGING THE DATA GAP: LEVERAGING AI TO ADDRESS DATA SCARCITY IN MEDICAL IMAGING

SN: 24076607

ABSTRACT

This section provides a brief overview of the methodology/results presented in the report.¹

Index Terms— One, two, three, four, five

1. APPLICATION DOMAIN AND CHALLENGES

1.1. The Application Domain: Medical Imaging

Medical imaging is pivotal in modern diagnostics, providing non-invasive methods to visualize internal structures of the human body. It encompasses a wide range of imaging modalities, including magnetic resonance imaging (MRI), computed tomography (CT), and X-rays, each offering different contrast mechanisms and spatial resolutions to support clinical diagnosis and treatment planning.

In recent years, medical imaging has been revolutionized by the advances in computer vision and deep learning technologies[1]. These technologies have significantly enhanced the capability of automatic image analysis systems, particularly in tasks like medical image segmentation (MedSeg), which involves partitioning images into regions corresponding to specific organs or lesions such as the chest, brain, abdomen, eye, and heart[2].

Among the imaging modalities, CT scans, MRI, and X-rays are especially critical in detecting and monitoring various health conditions. For instance, CT imaging has been extensively used in diagnosing lung infections during the COVID-19 pandemic, as evidenced by the segmentation of lung infection regions in CT scans[3, 4, 5].

The segmentation process itself plays a fundamental role in medical image analysis by enabling pixel-level identification of anatomical structures, facilitating precise diagnosis and personalized treatment[6]. Manual segmentation by radiologists, although accurate, is labor-intensive, time-consuming, and costly, thus driving the demand for automated methods that leverage deep learning[7].

Overall, medical imaging, supported by modern computational methods, forms the cornerstone of contempo-

rary healthcare, enabling earlier detection, better disease monitoring, and improved patient outcomes.

1.2. Current Challenges: Data Scarcity and Its Implications

Medical image segmentation has witnessed tremendous progress with deep learning, yet its success remains heavily dependent on the availability of large-scale, high-quality annotated datasets. Unfortunately, several persistent challenges hinder the widespread deployment of AI models in medical imaging, especially in real-world clinical environments.

1.2.1. Limited Annotated Datasets

Obtaining labeled medical imaging data is a labor-intensive and costly process, typically requiring the expertise of trained radiologists and high-end equipment. Manual annotation, such as pixel-wise segmentation, is especially time-consuming. As a result, the availability of large annotated datasets remains limited. Moreover, privacy regulations and patient confidentiality further restrict data sharing and public availability[7].

1.2.2. Bias and Generalizability Issues

Even when labeled datasets are available, they are often limited in diversity, both demographically and technically (e.g., variation in scanner models or acquisition protocols). This lack of heterogeneity leads to significant distribution shift problems when models trained on one dataset are deployed on another, ultimately affecting their generalization performance across populations and institutions.

1.2.3. Resource Constraints

Healthcare systems in low- and middle-income regions face acute shortages in data collection infrastructure and medical imaging resources. The high costs of annotation and hardware requirements for data processing place an additional burden on AI development in such contexts. As noted in the literature, while deep learning architectures like U-Nets are widely adopted in academic re-

¹The Blog is provided https://yushiran.github.io/ELEC0139_BLOG_SN24076607/ and GitHub project: https://github.com/yushiran/ELEC0139_BLOG_SN24076607

search, deploying them in under-resourced settings remains a formidable challenge.

1.3. The Case for AI/ML Technologies

This subsection makes the case for adopting machine learning and artificial intelligence technologies to address the challenges and improve outcomes in the application domain.

1.3.1. Efficient Data Utilization

AI techniques, especially deep learning architectures like U-Nets, are capable of learning meaningful spatial and semantic patterns even from limited labeled data. When designed appropriately, such models can achieve strong segmentation performance despite inherent challenges like noise and distribution shift [8]. Moreover, hybrid architectures such as U-Net++ and attention mechanisms have further improved efficiency and robustness [9].

1.3.2. Synthetic Data Generation

Synthetic data generation using generative adversarial networks (GANs) is one of the most promising avenues for alleviating labeled data scarcity. For instance, models like Cycle-GANs and conditional GANs have been used to generate high-fidelity synthetic MRI and ultrasound images that closely resemble real samples, including segmentation labels. These synthetic datasets, when used in model training, have shown comparable performance to real data [10, 11].

1.3.3. Self-Supervised Learning

Self-supervised learning (SSL) has gained traction in medical image analysis due to its ability to leverage vast amounts of unlabeled data. Techniques such as context restoration, multi-modal feature fusion, and attention-based pseudo labeling enable models to learn robust representations without requiring ground-truth masks [12, 13].

2. AI/ML SOLUTIONS TO DATA SCARCITY IN MEDICAL IMAGING

2.1. Self-Supervised Learning (SSL)

Self-supervised learning (SSL) enables the use of large amounts of unlabeled data to pretrain neural networks by defining pretext tasks—artificial supervision signals derived from the data itself. In medical imaging, this is particularly valuable, as obtaining labeled data is expensive and requires expert input.

Chen et al. (2019) [14] proposed a context restoration strategy tailored to the characteristics of medical

images. The method corrupts the spatial arrangement of an image by swapping randomly selected patches and then trains a convolutional neural network (CNN) to restore the original image. This process forces the network to learn semantic-level image representations, which are transferable to downstream tasks such as classification, localization, and segmentation.

2.1.1. Key Features of Context Restoration SSL

- **Semantic Feature Learning:** The network learns to recognize and correct structural inconsistencies, resulting in rich semantic representations.
- **Transferability:** Features learned through this method can initialize both encoder and decoder parts of downstream CNNs, especially useful for segmentation tasks where full image-to-image mapping is needed.
- **Simplicity:** The approach is simple to implement, with minimal changes to existing architectures and training pipelines.

2.1.2. Methodology

Let $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ be a set of unlabeled medical images. A corruption function \mathcal{R} generates a disordered image \tilde{x}_i :

$$\tilde{x}_i = \mathcal{R}(x_i)$$

A CNN model $g(\cdot)$ is then trained to restore the original image:

$$x_i = g(\tilde{x}_i) \approx f^{-1}(\tilde{x}_i)$$

The training objective is to minimize the pixel-wise L2 reconstruction loss:

$$\mathcal{L}_{\text{SSL}} = \|x_i - g(\tilde{x}_i)\|_2^2$$

The corruption function \mathcal{R} randomly selects and swaps image patches:

Algorithm 1 Image Context Disordering

- 1: Input: original image x_i
 - 2: Output: image with disordered context \tilde{x}_i
 - 3: for $t = 1$ to T do
 - 4: randomly select patch $p_1 \in x_i$
 - 5: randomly select patch $p_2 \in x_i$
 - 6: if $p_1 \cap p_2 = \emptyset$ then
 - 7: swap p_1 and p_2
 - 8: end if
 - 9: end for
-

The CNN model $g(\cdot)$ has two parts, as shown in Figure 1:

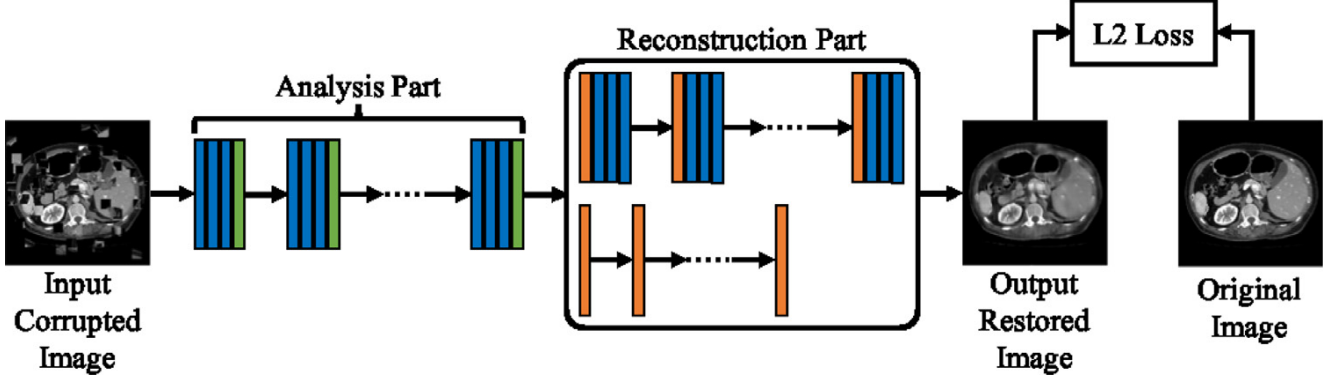


Fig. 1. General CNN architecture for context restoration SSL. Blue, green, and orange strides represent convolutional, downsampling, and upsampling units, respectively.

- **Analysis Part:** an encoder that extracts features from the disordered image. It may include convolutional layers, residual blocks [15], or inception modules [16].
- **Reconstruction Part:** a decoder that upsamples the features and reconstructs the image in correct spatial order.
- **Localization:** For abdominal organ localization in CT images, models initialized via context restoration outperformed those trained with auto-encoders or relative position tasks, especially under data-limited settings.
- **Segmentation:** In brain tumor segmentation using multi-modal MRI, models with context restoration pretraining achieved higher Dice scores and lower Hausdorff distances than all other SSL and baseline methods.

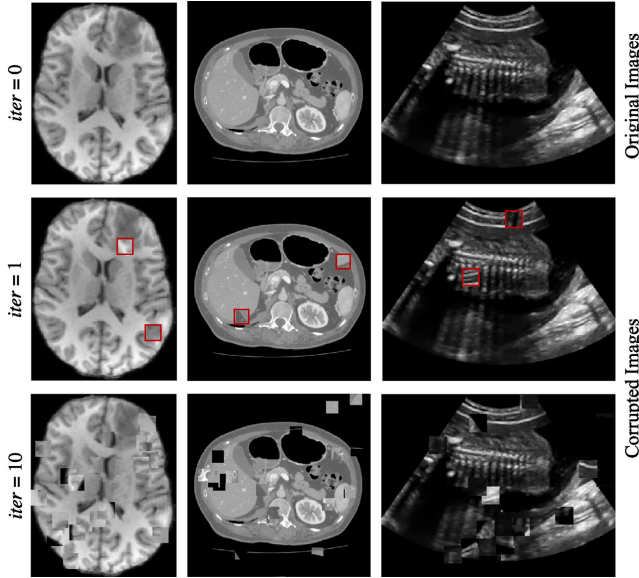


Fig. 2. Examples of training images for self-supervised context disordering. The second column highlights swapped patches after the first iteration.

2.1.3. Applications and Evaluation

- **Classification:** On fetal ultrasound images, context restoration pretraining improved the F1-score by over 7 percentage points compared to random initialization with only 25% of training data.

2.1.4. Benefits

- Reduces dependence on labeled data by leveraging vast pools of unlabeled medical images.
- Improves model performance under limited supervision conditions, particularly in small-sample settings.
- Generalizes well across modalities (ultrasound, CT, MRI) and tasks (classification, localization, segmentation).

2.2. Reinforcement Learning (RL)

Reinforcement Learning (RL) is a powerful machine learning paradigm in which an agent learns to interact with its environment by receiving feedback in the form of rewards. Unlike supervised learning, which relies heavily on large-scale annotated datasets, RL can operate effectively with minimal labeled data, making it particularly attractive in medical imaging domains where data scarcity is a major challenge [17].

An RL framework is typically defined by a set of core components: state (the environment observation), action (possible moves the agent can make), reward (feedback signal guiding learning), and policy (the decision-

making strategy). Depending on whether the environment is explicitly modeled, RL approaches are broadly categorized into model-free and model-based methods. Model-free methods, such as DQN[18] and A2C[19, 20], learn policies directly through interaction, while model-based approaches attempt to learn a transition model to improve sample efficiency—particularly important in low-data regimes.

2.2.1. Mathematical Formulation of Reinforcement Learning

Reinforcement learning problems are often modeled as a Markov Decision Process (MDP), defined by a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$, where:

- \mathcal{S} is the set of possible states,
- \mathcal{A} is the set of actions,
- $\mathcal{P}(s'|s, a)$ is the transition probability function,
- $\mathcal{R}(s, a)$ is the reward received after taking action a in state s ,
- $\gamma \in [0, 1]$ is the discount factor for future rewards.

The goal is to learn a policy $\pi(a|s)$ that maximizes the expected cumulative reward:

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

The value function for a state under policy π is:

$$V^{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right]$$

The action-value function (Q-function) is:

$$Q^{\pi}(s, a) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right]$$

An optimal policy π^* satisfies:

$$Q^{\pi^*}(s, a) = \max_{\pi} Q^{\pi}(s, a)$$

This formulation allows reinforcement learning agents to learn optimal decision-making strategies through trial-and-error, without requiring large-scale labeled data. This is especially beneficial in medical imaging applications such as classification, registration, or synthesis, where annotated datasets are scarce.

The generic reinforcement learning procedure outlined in Algorithm 2 provides a high-level framework for training RL agents. It begins with initializing the policy parameters and iteratively updates them based on

the agent’s interactions with the environment. At each step, the agent selects an action according to its current policy, observes the resulting reward and next state, and updates the policy parameters to improve performance. This iterative process continues until the policy converges to an optimal or near-optimal solution.

Algorithm 2 highlights the flexibility of RL in handling diverse tasks, as it does not rely on predefined labels but instead learns directly from the environment’s feedback. This makes it particularly suitable for medical imaging applications, where labeled data is often scarce or expensive to obtain.

Algorithm 2 Generic Reinforcement Learning Procedure

- 1: Input: Environment \mathcal{E} , initial policy π_{θ}
 - 2: Initialize policy parameters θ randomly
 - 3: for each episode do
 - 4: Initialize state s_0
 - 5: for each step $t = 0, 1, 2, \dots$ until terminal state do
 - 6: Select action $a_t \sim \pi_{\theta}(a_t|s_t)$
 - 7: Execute a_t , observe reward r_t and next state s_{t+1}
 - 8: Update policy parameters θ using transition (s_t, a_t, r_t, s_{t+1})
 - 9: end for
 - 10: end for
 - 11: Output: Trained policy π_{θ}
-

2.2.2. Applications of Reinforcement Learning in Medical Imaging

RL has been successfully applied to a wide range of medical imaging tasks, including image classification, landmark localization, lesion detection, segmentation, image registration, and radiotherapy planning. These applications span multiple anatomical sites (e.g., brain, lung, prostate) and imaging modalities (e.g., MRI, CT, ultrasound), as summarized in Figure 3.

Importantly, RL offers several key mechanisms to alleviate data scarcity in medical imaging:

- **Minimal dependence on annotations:** RL agents can learn optimal behaviors by interacting with environments, reducing reliance on large-scale annotated datasets.
- **Higher sample efficiency:** Especially in model-based RL, agents require fewer interactions to achieve comparable performance, making them well-suited for small datasets.
- **Active data selection:** RL-based frameworks have been proposed to select the most informative samples for annotation or training, optimizing the use of limited labeled data.

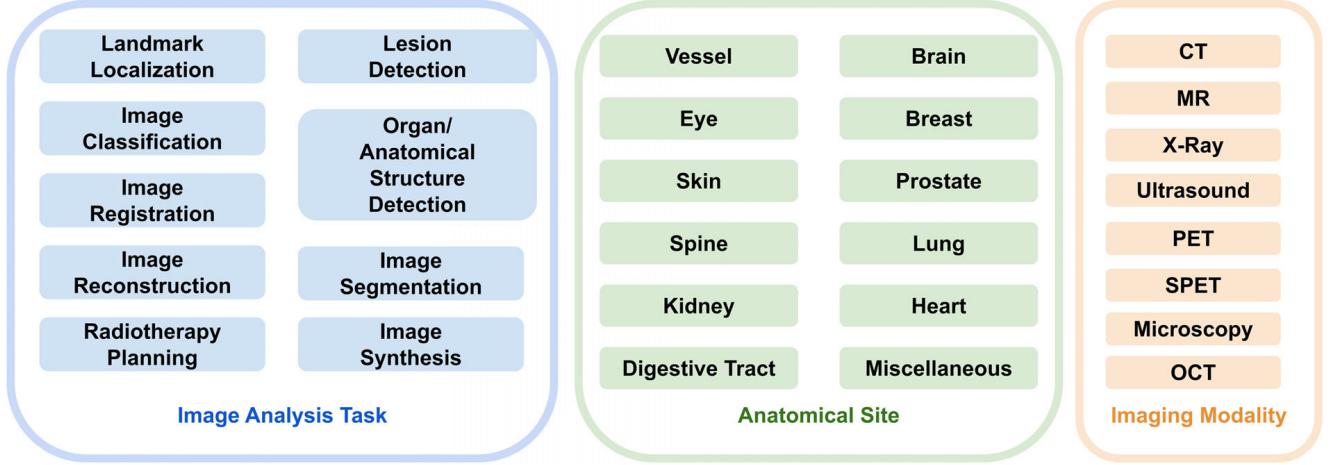


Fig. 3. Blue box covers image analysis tasks; green box covers anatomical sites; yellow box covers imaging modalities.

- Combination with generative models: RL can be integrated with GANs or VAEs to select high-quality synthetic samples for augmentation, effectively enhancing dataset diversity.

Overall, reinforcement learning not only reduces the burden of manual annotation but also promotes the development of data-efficient, adaptive, and goal-driven medical image analysis systems. Its ability to model complex sequential decision-making makes it a promising tool for next-generation clinical AI.

2.3. Generative Models for Medical Image Synthesis

The scarcity of large, diverse, and annotated medical imaging datasets remains a significant bottleneck in developing robust AI models for clinical applications. Generative models, particularly Generative Adversarial Networks (GANs)[21] and diffusion models[22], have emerged as powerful tools to mitigate this challenge by synthesizing high-fidelity medical images that augment limited real-world datasets [23]. These models learn the underlying data distribution of training images and generate novel samples that preserve anatomical and pathological features while ensuring patient privacy.

2.3.1. Generative Adversarial Networks (GANs)

GANs consist of a generator (creates synthetic images) and a discriminator (distinguishes real from synthetic images). Through adversarial training, the generator improves realism[21].

For example, Upadhyay et al. (2024)[24] proposed a GAN-based framework to generate synthetic lung lesions mimicking ground glass nodules (GGNs), addressing the data scarcity issue in computer-aided diagnosis

systems. The model consists of a generator and a discriminator trained adversarially to produce realistic synthetic GGNs.

The generator employs a U-Net-like architecture to synthesize GGNs[25], while the discriminator uses convolutional layers[15] to distinguish real from synthetic images. The loss function combines adversarial loss with pixel-wise reconstruction loss to ensure both realism and anatomical accuracy.

The model consists of three key components:

- Generator (G): SRGAN-based network that synthesizes pulmonary nodules from masked input images
- ROI Discriminator (D_{ROI}): ResNet-based classifier operating on nodule regions (red path in Fig.4)
- Whole Image Discriminator (D_{whole}): Parallel ResNet evaluating full contextual realism (blue path)

The composite loss combines adversarial and similarity terms for both discriminators:

$$\mathcal{L}_{DSRGAN} = (\mathcal{L}_{sim} + \mathcal{L}_{adv})_{whole} + (\mathcal{L}_{sim} + \mathcal{L}_{adv})_{ROI} \quad (1)$$

$$\mathcal{L}_{adv} = \sum_{n=1}^N -\log D(G(x)) \quad (2)$$

$$\mathcal{L}_{sim}(x, y) = 1 - \frac{(2\mu_x\mu_y + C_1) + (\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad (3)$$

where μ, σ denote mean/variance of image patches, C_1, C_2 stabilize division.

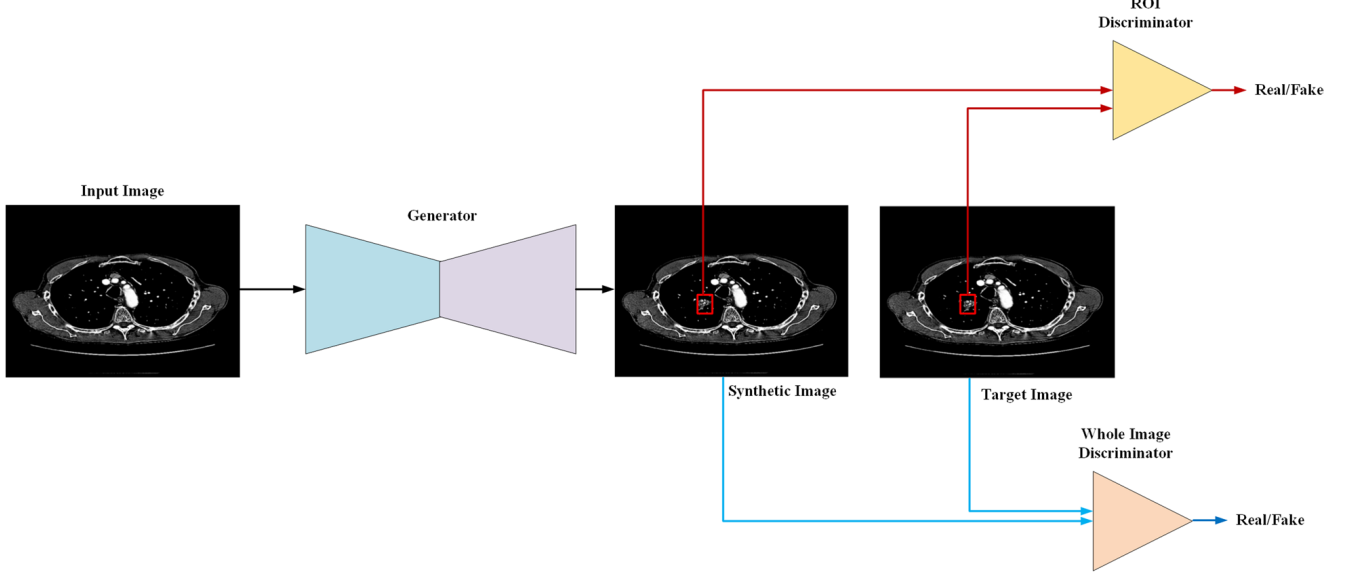


Fig. 4. A schematic overview of a diffusion model in training and sampling settings. In the top row, the diffusion model is trained and creates a Markov chain to add Gaussian noise to the real images, resulting in a noise vector z . The model then reverses the Markov chain by predicting the next state of the image from the current noisy state, which is equivalent to denoising the image. During sampling (bottom row), the model can generate synthetic images by starting from a random noise vector and applying the reverse Markov chain.

The result of the GAN training is a generator capable of producing synthetic GGNs that closely resemble real lesions, as shown in Figure 5. The generated images can be used to augment existing datasets, improving the performance of downstream tasks such as classification and segmentation.

$$q_t(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_0, (1 - \alpha_t)\mathbf{I}), \quad \alpha_t = \prod_{s=1}^t (1 - \beta_s) \quad (5)$$

The perturbed data distribution is defined as:

$$p_{\alpha_t}(\tilde{x}) = \int p_{data}(x)q_{\alpha_t}(\tilde{x}|x)dx \quad (6)$$

with noise scales chosen such that $X_N \approx \mathcal{N}(0, \mathbf{I})$.

The variational Markov chain in the reverse direction is parameterized as:

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}\left(x_{t-1}; \frac{1}{\sqrt{1 - \beta_t}}(x_t + \beta_t s_{\theta}(x_t, t)), \beta_t \mathbf{I}\right) \quad (7)$$

The model is trained with a re-weighted ELBO variant:

$$\theta^* = \arg \min_{\theta} \sum_{t=1}^N (1 - \alpha_t) \mathbb{E}_{p_{data}(x)} \mathbb{E}_{q_t(\tilde{x}|x)} \|s_{\theta}(\tilde{x}, t) - \nabla_x \log q_t(\tilde{x}|x)\|_2^2 \quad (8)$$

After obtaining the optimal model s_{θ}^* , samples's generation is as shown in Algorithm 3.

2.3.2. Diffusion Models

Diffusion models are a class of generative models that learn to generate data by reversing a diffusion process. They have gained popularity due to their ability to produce high-quality samples and have been successfully applied in various domains, including image synthesis, text generation, and audio processing[26].

Figure 6 illustrates the training and sampling process of the diffusion model, showcasing how noise is added and subsequently removed to generate synthetic images.

Consider a sequence of positive noise scales $0 < \beta_1, \dots, \beta_N < 1$. For each training data point $x_0 \sim p_{data}(x)$, construct a discrete Markov chain $\{X_0, X_1, \dots, X_N\}$ where:

$$p(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}) \quad (4)$$

The marginal distribution after t steps becomes:

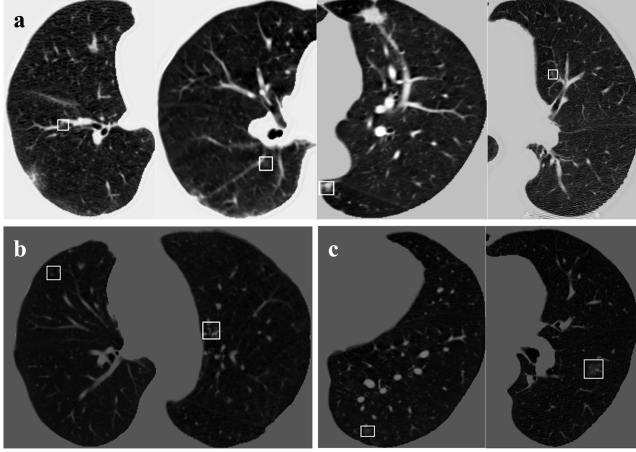


Fig. 5. Examples of synthetic ground glass nodules (GGNs), the GGNs were categorised by physicians to four categories: confidently fake, leaning fake, leaning real, and confidently real. a Synthetic GGNs classified as “real” by clinicians. b Synthetic GGNs with less convincing generated lesions (classified as “leaning fake”). c A real GGNs in the original LIDC-IDRI dataset

Algorithm 3 DDPM Ancestral Sampling

```

1: Initialize  $x_N \sim \mathcal{N}(0, \mathbf{I})$ 
2: for  $t = N$  downto 1 do
3:    $x_{t-1} = \frac{1}{\sqrt{1-\beta_t}}(x_t + \beta_t s_\theta^*(x_t, t)) + \sqrt{\beta_t} z_t, z_t \sim \mathcal{N}(0, \mathbf{I})$ 
4: end for
5: Return  $x_0$ 

```

Diffusion models are particularly well-suited for medical image generation due to their ability to produce high-quality, diverse, and anatomically accurate synthetic images. Their iterative denoising process ensures fine-grained control over the generated data, preserving critical medical details. Additionally, diffusion models are robust to noise and can effectively model complex data distributions, making them ideal for handling the variability and precision required in medical imaging. These characteristics make diffusion models a powerful tool for augmenting datasets, improving model generalization, and addressing data scarcity challenges in medical imaging applications.

3. ETHICAL CONSIDERATIONS IN APPLYING AI TO MEDICAL IMAGING

The integration of AI into medical imaging has yielded impressive results, but also raises several ethical concerns. These issues must be thoroughly addressed to ensure safe, equitable, and trustworthy deployment of AI systems in healthcare.

3.1. Data Privacy and Security

The collection and use of medical imaging data raise significant privacy concerns due to the personal nature of health information. Patient data is strictly protected by various regulations such as HIPAA in the United States, GDPR in Europe, and similar frameworks globally. These regulations impose strict requirements on data sharing and access, which can significantly impede the development of robust AI models [23].

Synthetic data generation techniques, as discussed in Section 2, offer a promising solution to these privacy challenges. By creating artificial medical images that preserve the statistical properties and clinical relevance of real images without containing actual patient information, these approaches can effectively circumvent many privacy concerns:

- **Risk Mitigation:** Synthetic data eliminates the risk of exposing protected health information (PHI), as the generated images do not correspond to real patients. This significantly reduces the regulatory burden and potential legal liability associated with data breaches.
- **Enhanced Data Sharing:** Synthetic datasets can be more freely shared across institutions and international boundaries, facilitating collaborative research and development without privacy impediments.
- **Data Augmentation:** As demonstrated in our GAN and diffusion model implementations, synthetic images can augment limited real datasets, addressing both privacy concerns and data scarcity simultaneously.

However, synthetic data is not without its own security considerations. Particularly, models like GANs might inadvertently memorize training examples, potentially leading to data leakage if not properly safeguarded. Additionally, adversarial attacks on these generative models could potentially extract sensitive information from the training data. Rigorous security measures, including proper model evaluation for memorization, differential privacy techniques during training, and robust access controls, must be implemented to ensure synthetic data approaches maintain strong privacy guarantees [23].

The self-supervised learning approaches described earlier provide another privacy-preserving advantage: they can extract valuable information from unlabeled data without requiring detailed annotations that might contain sensitive information. By learning from the inherent structure of images rather than explicit labels, techniques like context restoration minimize exposure to privacy-sensitive metadata.

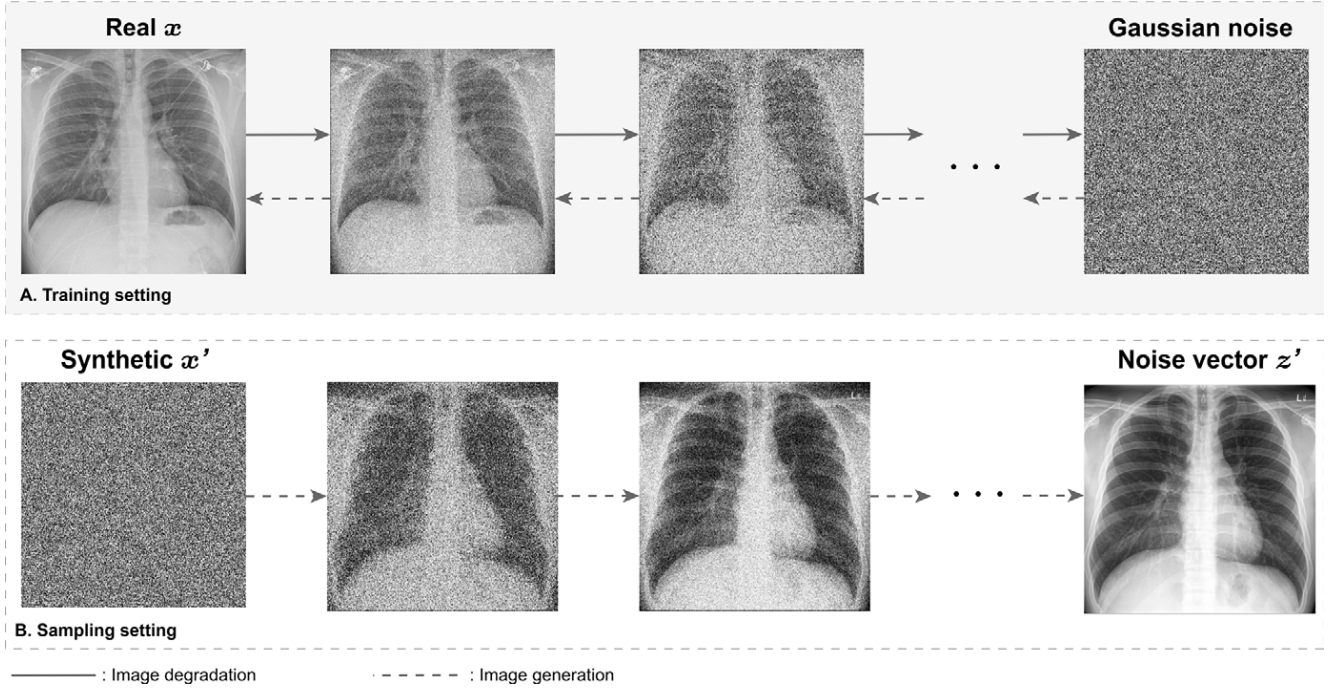


Fig. 6. A schematic overview of a diffusion model in training and sampling settings. In the top row, the diffusion model is trained and creates a Markov chain to add Gaussian noise to the real images, resulting in a noise vector z' . The model then reverses the Markov chain by predicting the next state of the image from the current noisy state, which is equivalent to denoising the image. During sampling (bottom row), the model can generate synthetic images by starting from a random noise vector and applying the reverse Markov chain.

3.2. Bias and Fairness

AI systems for medical imaging are inherently shaped by the data used to train them, making them susceptible to perpetuating or even amplifying existing biases in healthcare. This is particularly concerning in the context of data scarcity, where models might be developed using imbalanced or non-representative datasets [23].

Several types of bias can manifest in medical imaging AI:

- **Demographic Bias:** When training data lacks diversity across age, gender, ethnicity, or socioeconomic factors, resulting models may perform disproportionately poorly on underrepresented groups. For example, models trained predominantly on data from certain ethnic populations may show reduced accuracy when applied to different demographic groups.
- **Technical Bias:** Variations in imaging equipment, acquisition protocols, and institutional practices introduce substantial heterogeneity in medical images. Models trained on data from high-end scanners may perform poorly on images from lower-resource settings, potentially exacerbating

healthcare disparities.

- **Selection Bias:** The process of collecting training data often introduces sampling biases. For instance, data from academic medical centers may overrepresent rare or complex cases compared to community hospitals.

The generative approaches discussed in Section 2, while addressing data scarcity, introduce their own fairness considerations. GANs and diffusion models tend to capture and potentially amplify patterns present in their training data. If the training data contains biases, synthetic data generated from these models may inherently encode and propagate these biases, potentially worsening the problem rather than solving it.

To mitigate these risks, synthetic data generation should be specifically designed with fairness in mind:

- **Balanced Data Generation:** Generative models can be explicitly conditioned to produce balanced distributions across demographic factors or modalities, potentially oversampling underrepresented groups.
- **Fairness-Aware Training:** Incorporating fairness constraints or adversarial debiasing techniques

into generative model training can help reduce the transfer of biases to synthetic data.

- **Diverse Data Sources:** Incorporating data from diverse sites and populations, even if in small quantities, can help generative models capture broader variations in anatomical structures and pathologies.

3.3. Transparency and Explainability

The "black-box" nature of many advanced AI systems used in medical imaging poses significant challenges for clinical adoption and regulatory approval. Complex models like deep generative networks and self-supervised learning approaches often lack transparency in their decision-making processes, making it difficult for healthcare professionals to understand and trust their outputs[8].

In medical contexts, where decisions directly impact patient outcomes, this lack of explainability is particularly problematic:

- **End-to-end Generative Models:** GANs and diffusion models operate as black boxes, taking inputs and producing synthetic images through complex transformations that are not easily interpretable. The multi-layer, non-linear nature of these models makes it virtually impossible to trace exactly how specific features in the generated images were constructed.
- **Self-Supervised Learning:** While SSL methods effectively learn from unlabeled data, the representations they develop are often abstract and difficult to map to clinically meaningful features. The pretext tasks (like context restoration) may have little direct relationship to the downstream diagnostic tasks.
- **Reinforcement Learning:** Among the approaches discussed, RL may offer slightly better explainability through its explicit reward functions and state-action mappings, but complex neural network policies still suffer from opacity in their internal reasoning.

To address these challenges, several approaches are being developed to enhance the explainability of AI in medical imaging:

- **Counterfactual Explanations:** Generating "what-if" scenarios that demonstrate how changes to the input would affect the output helps users understand the model's decision boundaries.

- **Layer-wise Relevance Propagation:** This technique decomposes predictions into contributions from individual input features, creating heatmaps that visualize important regions.
- **Feature Disentanglement:** Particularly for generative models, encouraging the separation of clinically relevant features (e.g., anatomical structures, pathologies) into interpretable latent dimensions improves transparency.

3.4. Accountability and Governance

In medical imaging, where AI impacts patient care, robust accountability and governance frameworks are critical. The discussed AI approaches—generative models, self-supervised learning, and reinforcement learning—pose unique challenges, especially in data-scarce contexts.

- **Clinical Validation:** AI systems using limited or synthetic data must undergo rigorous testing across diverse populations with clear performance metrics and post-deployment monitoring.
- **Responsibility Assignment:** Clear accountability frameworks are needed to define responsibility for errors, especially with models trained on synthetic or limited data.
- **Regulatory Oversight:** Evolving frameworks, like those from the FDA, must address challenges specific to data-scarce and generative AI systems.

Governance of synthetic data requires attention to:

- **Quality Control:** Ensure synthetic images meet clinical standards and represent pathological features accurately.
- **Provenance Tracking:** Maintain records distinguishing real and synthetic data for transparency.
- **Continuous Evaluation:** Regularly reassess models trained on synthetic data to ensure reliability.

For self-supervised and reinforcement learning:

- **Pretext Task Validation:** Ensure self-supervised tasks produce clinically relevant representations.
- **Reward Function Oversight:** Design RL reward functions collaboratively with clinicians to ensure meaningful outcomes.
- **Update Protocols:** Define clear guidelines for model updates and recertification.

Multidisciplinary committees and international standards are essential to ensure safety, efficacy, and equity in deploying AI in data-scarce medical domains.

4. CONCLUSION

This paper explored how artificial intelligence and machine learning technologies can address data scarcity challenges in medical imaging. We examined three key approaches: self-supervised learning, reinforcement learning, and generative models. Self-supervised learning effectively leverages unlabeled data through pretext tasks like context restoration. Reinforcement learning offers learning from limited feedback rather than extensive labeled datasets. Generative models, including GANs and diffusion models, synthesize realistic medical images to augment limited datasets while preserving patient privacy.

While these technologies show great promise, ethical considerations around privacy, bias, transparency, and governance remain crucial. With continued research and responsible implementation, AI technologies can help overcome data limitations, potentially improving healthcare access and outcomes across diverse clinical settings.

5. REFERENCES

- [1] Ashwini Kumar Upadhyay and Ashish Kumar Bhandari, "Advances in Deep Learning Models for Resolving Medical Image Segmentation Data Scarcity Problem: A Topical Review," *Archives of Computational Methods in Engineering*, vol. 31, no. 3, pp. 1701–1719, Apr. 2024.
- [2] Sushu Sushanki, Ashish Kumar Bhandari, and Amit Kumar Singh, "A Review on Computational Methods for Breast Cancer Detection in Ultrasound Images Using Multi-Image Modalities," *Archives of Computational Methods in Engineering*, vol. 31, no. 3, pp. 1277–1296, Apr. 2024.
- [3] Liangliang Liu, Jianhong Cheng, Quan Quan, Fang-Xiang Wu, Yu-Ping Wang, and Jianxin Wang, "A survey on U-shaped networks in medical image segmentations," *Neurocomputing*, vol. 409, pp. 244–258, 2020.
- [4] Nahian Siddique, Sidike Paheding, Colin P. Elkin, and Vijay Devabhaktuni, "U-net and its variants for medical image segmentation: A review of theory and applications," *IEEE Access*, vol. 9, pp. 82031–82057, 2021.
- [5] Xin Yi, Ekta Walia, and Paul Babyn, "Generative adversarial network in medical imaging: A review," *Medical Image Analysis*, vol. 58, pp. 101552, 2019.
- [6] Sonu Kumar, Ashish Kumar Bhandari, Aditya Raj, and Kirti Swaraj, "Triple Clipped Histogram-Based Medical Image Enhancement Using Spatial Frequency," *IEEE transactions on nanobioscience*, vol. 20, no. 3, pp. 278–286, July 2021.
- [7] Nima Tajbakhsh, Laura Jeyaseelan, Qian Li, Jeffrey N. Chiang, Zhihao Wu, and Xiaowei Ding, "Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation," *Medical Image Analysis*, vol. 63, pp. 101693, 2020.
- [8] Poonam Rani Verma and Ashish Kumar Bhandari, "Role of Deep Learning in Classification of Brain MRI Images for Prediction of Disorders: A Survey of Emerging Trends," *Archives of Computational Methods in Engineering*, vol. 30, no. 8, pp. 4931–4957, Nov. 2023.
- [9] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang, "UNet++: A Nested U-Net Architecture for Medical Image Segmentation," in *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, Danail Stoyanov,

- Zeike Taylor, Gustavo Carneiro, Tanveer Syeda-Mahmood, Anne Martel, Lena Maier-Hein, João Manuel R.S. Tavares, Andrew Bradley, João Paulo Papa, Vasileios Belagiannis, Jacinto C. Nascimento, Zhi Lu, Sailesh Conjeti, Mehdi Moradi, Hayit Greenspan, and Anant Madabhushi, Eds., Cham, 2018, pp. 3–11, Springer International Publishing.
- [10] Andrew Gilbert, Maciej Marciniak, Cristobal Roderio, Pablo Lamata, Eigil Samset, and Kristin Mcleod, “Generating synthetic labeled data from existing anatomical models: An example with echocardiography segmentation,” *IEEE Transactions on Medical Imaging*, vol. 40, no. 10, pp. 2783–2794, 2021.
 - [11] Hoo-Chang Shin, Neil A. Tenenholz, Jameson K. Rogers, Christopher G. Schwarz, Matthew L. Senjem, Jeffrey L. Gunter, Katherine P. Andriole, and Mark Michalski, “Medical Image Synthesis for Data Augmentation and Anonymization Using Generative Adversarial Networks,” in *Simulation and Synthesis in Medical Imaging*, Ali Gooya, Orcun Goksel, Ipek Oguz, and Ninon Burgos, Eds., Cham, 2018, pp. 1–11, Springer International Publishing.
 - [12] Krishna Chaitanya, Neerav Karani, Christian F. Baumgartner, Ertunc Erdil, Anton Becker, Olivio Donati, and Ender Konukoglu, “Semi-supervised task-driven data augmentation for medical image segmentation,” *Medical Image Analysis*, vol. 68, pp. 101934, 2021.
 - [13] Hao Zheng, Jun Han, Hongxiao Wang, Lin Yang, Zhuo Zhao, Chaoli Wang, and Danny Z. Chen, “Hierarchical Self-supervised Learning for Medical Image Segmentation Based on Multi-domain Data Aggregation,” in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2021*, Marleen de Bruijne, Philippe C. Cattin, Stéphane Cotin, Nicolas Padoy, Stefanie Speidel, Yefeng Zheng, and Caroline Essert, Eds., Cham, 2021, pp. 622–632, Springer International Publishing.
 - [14] Liang Chen, Paul Bentley, Kensaku Mori, Kazunari Misawa, Michitaka Fujiwara, and Daniel Rueckert, “Self-supervised learning for medical image analysis using image context restoration,” *Medical Image Analysis*, vol. 58, pp. 101539, Dec. 2019.
 - [15] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.
 - [16] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna, “Rethinking the inception architecture for computer vision,” *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2818–2826, 2015.
 - [17] Mingzhe Hu, Jiahan Zhang, Luke Matkovic, Tian Liu, and Xiaofeng Yang, “Reinforcement learning in medical image analysis: Concepts, applications, challenges, and future directions,” *Journal of Applied Clinical Medical Physics*, vol. 24, no. 2, pp. e13898, 2023.
 - [18] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
 - [19] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov, “Proximal policy optimization algorithms,” *ArXiv*, vol. abs/1707.06347, 2017.
 - [20] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in *International Conference on Machine Learning*, 2016.
 - [21] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio, “Generative adversarial nets,” in *Neural Information Processing Systems*, 2014.
 - [22] Jascha Narain Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics,” *ArXiv*, vol. abs/1503.03585, 2015.
 - [23] Lennart R. Koetzier, Jie Wu, Domenico Mastrodicasa, Aline Lutz, Matthew Chung, W. Adam Koszek, Jayanth Pratap, Akshay S. Chaudhari, Pranav Rajpurkar, Matthew P. Lungren, and Martin J. Willeminck, “Generating Synthetic Data for Medical Imaging,” *Radiology*, vol. 312, no. 3, pp. e232471, Sept. 2024.
 - [24] Zhixiang Wang, Zhen Zhang, Ying Feng, Lizza E. L. Hendriks, Razvan L. Miclea, Hester Gietema, Janna Schoenmaekers, Andre Dekker, Leonard Wee, and Alberto Traverso, “Generation of synthetic ground glass nodules using generative adversarial networks

(GANs),” *European Radiology Experimental*, vol. 6, no. 1, pp. 59, Nov. 2022.

- [25] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi, “Photo-realistic single image super-resolution using a generative adversarial network,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 105–114.
- [26] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole, “Score-Based Generative Modeling through Stochastic Differential Equations,” Feb. 2021.